

# An Adaptable Binary Entropy Coder

Aaron Kiely and Matthew Klimesh

Jet Propulsion Laboratory, California Institute of Technology  
4800 Oak Grove Drive, Mail Stop 238-420, Pasadena, CA 91109  
e-mail: {aaron, klimesh}@shannon.jpl.nasa.gov

**Abstract**— We present a novel entropy coding technique which is based on recursive interleaving of variable-to-variable length binary source codes. The encoding is adaptable in that each bit to be encoded may have an associated probability estimate which depends on previously encoded bits. The technique can achieve arbitrarily small redundancy, and may have advantages over arithmetic coding, including most notably the admission of a simple and fast decoder. We discuss code design and performance estimation methods, as well as practical encoding and decoding algorithms.

## 1 Introduction

In data compression algorithms the need frequently arises to compress a binary sequence in which each bit has some estimated distribution, i.e., probability of being equal to zero. If long runs of bits have nearly identical distributions, then simple source codes, most notably Golomb's runlength codes [1], are quite efficient. However, in many practical situations, not only does the distribution vary from bit to bit, but it is desirable to have the estimated distribution for a bit depend on the values of earlier bits. Accommodating such a dynamically changing probability estimate is tricky because the decoder must make the same estimates as the encoder. Thus, before the  $i$ th bit can be decoded, the value of the first  $i - 1$  bits must be determined. This complication makes it difficult to efficiently use simple source codes such as runlength codes.

To our knowledge, currently the only efficient encoding methods in this case are arithmetic coding [2, 3] and the relatively unknown technique called interleaved entropy coding [4], which is a generalization of the "block Melcode" [5]. In this paper, we describe a new entropy coding technique which is a generalization of the interleaved entropy coding method. The technique efficiently encodes a binary source with a bit-wise adaptive probability estimate by recursively encoding groups of bits with similar distributions, ordering the output in a way that is suited to the decoder. As a result, the decoder has low complexity.

The functionality of our coding technique is essentially the same as that of binary arithmetic coding; however, our technique does not yield an arithmetic code and there are many practical differences. Arithmetic encoding of one bit requires a few arithmetic operations and, unless approximations are made, at least one multiplication. Our encoder requires no arithmetic operations except those needed to choose a

---

The work described was funded by the TMOD Technology Program and performed at the Jet Propulsion Laboratory, California Institute of Technology under contract with the National Aeronautics and Space Administration.

code index based on the bit distribution; however it requires some bookkeeping and bit manipulation operations. Our encoder requires more memory than arithmetic coding. Arithmetic decoders are generally of similar complexity to the encoders, but our decoder is much simpler than our encoder: it needs fewer operations than the encoder, and requires only a small amount of memory.

## 1.1 The Source Coding Problem

We examine the problem of compressing a sequence of bits  $b_1, b_2, \dots$  from a random source. The estimate of source probability  $p_i = \text{Prob}[b_i = 0]$  may depend on the values of the source sequence prior to index  $i$ , and on any other information available to both the encoder and decoder. This dependence encompasses both adaptive probability estimation as well as correlations or memory in the source. Consequently, efficient encoding requires a bit-wise adaptable encoder. We are not concerned here with methods of modeling the source, and so we make no distinction between the actual and estimated source distributions.

Without loss of generality, we will assume that  $p_i \geq 1/2$  for each index  $i$ . If this were not the case for some  $p_i$ , we could simply invert bit  $b_i$  before encoding to make it so (and this inversion can clearly be duplicated in the decoder).

We also assume that the decoder can determine when decoding is complete. In practice, this often occurs automatically, or straightforward methods can be used, such as transmitting the sequence length prior to the compressed sequence. Such methods will not be addressed here.

Although we only discuss the compression of binary sequences, given any nonbinary source we can assign prefix-free binary codewords to source symbols to produce a binary stream. Thus the technique can be applied to nonbinary sources as well.

## 1.2 The Recursive Interleaved Entropy Coding Concept

In this section we give an overview of how the entropy coding technique works and why it yields low redundancy. To simplify the explanation, some of the processing details are omitted until Section 2.

Since, by assumption, each bit has probability of zero at least  $1/2$ , we are concerned with the probability region  $[1/2, 1]$ . We partition this region into several narrow intervals, and with each interval we associate a bin that will be used to store bits. When bit  $b_i$  arrives, we place it into the bin corresponding to the interval containing  $p_i$ . Because each interval spans a small probability range, all of the bits in a given bin have nearly the same probability of being zero, and we can think of each bin as corresponding to some nominal probability value.

For each bin (except the leftmost bin, which contains probability  $1/2$ ) we specify an exhaustive prefix-free set of binary codewords. When the bits collected in a bin form one of these codewords, we delete these bits from the bin and encode the value of the codeword by placing one or more new bits in other bins<sup>1</sup>. This process is

---

<sup>1</sup>The ordering of the new bits in a bin is not straightforward and we save these details for

conveniently described using a binary tree. Each codeword is assigned to a terminal node in the tree, non-terminal nodes are labeled with a destination bin, and the branch labels (each a zero or one) correspond to the output bits that are placed in the destination bins.

For example, Figure 1 shows a tree that might be used for a bin with nominal probability 0.9. The prefix-free codeword set for this bin is  $\{00, 01, 1\}$ , shown as labels of the terminal nodes in the tree. If the codeword to be processed in the bin is 00, which occurs with probability approximately 0.81, we place a zero in the bin that contains probability 0.81. If the codeword is 1, first we place a one in the bin containing probability 0.81, which indicates that the codeword is something other than 00, then we place a zero in the bin containing probability 0.53 because, given that the codeword is not 00, the conditional probability that the codeword is 1 is approximately 0.53. We can see that this process might contribute to data compression because the most likely codeword is 00, which is represented using a single bit.

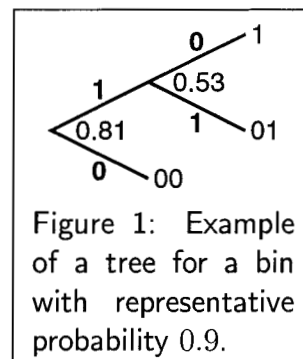


Figure 1: Example of a tree for a bin with representative probability 0.9.

For the leftmost bin we do not define a tree such as the one in Figure 1. Instead, bits in this bin form the encoder's output. Bits that reach the first bin have probability of being zero very close to  $1/2$  and are thus nearly incompressible, so transmitting these bits uncoded does not add much redundancy.

During the encoding process, bits arrive in various bins either directly from the source or as a result of processing codewords in other bins. Our goal is to have bits migrate to the leftmost (uncoded) bin, where they are transmitted. To accomplish this, we impose the constraint on the design of our trees that all new bits resulting from the processing of each codeword must be placed in bins strictly to the left of the bin in which the codeword was formed. Apart from our desire to move bits to the left, this constraint also prevents encoded information from traveling in "loops", which would make coding difficult or impossible. Thus if a bin has nominal probability  $p$ , we would like the probability of a zero for each output bit to be in the range  $[1/2, p)$ . A tree with this property is said to be *useful* at  $p$ . Perhaps surprisingly, useful trees exist everywhere:

**Theorem** *For any given probability value  $p \in (1/2, 1)$ , there exists a useful tree, i.e., one with the property that all output bits have probability of zero in the range  $[1/2, p)$ .*

This is proved by constructing an infinite family of trees for which at least one tree is useful at any given  $p \in (1/2, 1)$ . Figure 2 illustrates this construction. We omit the details of the proof.

When we reach the end of the bit sequence to be encoded and no codewords remain in any bin, there will generally be partially formed codewords in one or more bins. Since these bits are needed for decoding, we append one or more extra bits to each of these partial codewords to form complete codewords which are then processed in the normal manner<sup>2</sup>.

---

## Section 2.2.

<sup>2</sup>The method of selecting these extra bits that "flush" the encoder is relatively unimportant.

We can see that some redundancy is present in this system because the bins have positive width — the probability associated with a bit that arrives in a bin will usually not exactly equal the bin's nominal probability, and bits in the leftmost bin are transmitted uncompressed even though they may not have probability of zero exactly equal to  $1/2$ . However, we can reduce the redundancy by increasing the the number of bins and/or the size of the trees.

In practice, the encoder and decoder do not keep track of probability values. Instead, each bin is assigned an index, starting with index 1 corresponding to the leftmost (uncoded) bin. At each non-terminal node in the tree we identify the index, rather than the nominal probability value, of the bin to which the associated output bit is mapped. The constraint we impose on encoder design is that each output bit from the tree for bin  $j$  must be mapped to a bin with index strictly less than  $j$ . No computations involving probability values are needed apart from those which may be required to map each input bit  $b_i$  to the appropriate bin index.

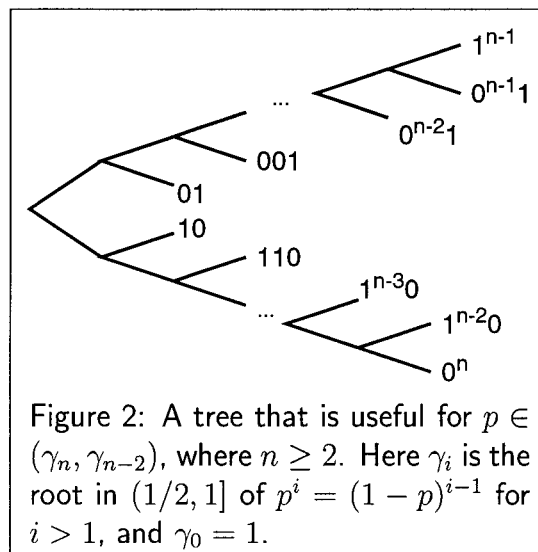


Figure 2: A tree that is useful for  $p \in (\gamma_n, \gamma_{n-2})$ , where  $n \geq 2$ . Here  $\gamma_i$  is the root in  $(1/2, 1]$  of  $p^i = (1 - p)^{i-1}$  for  $i > 1$ , and  $\gamma_0 = 1$ .

For example, a five bin encoder is defined by the trees shown in Figure 3. Figure 3(c) indicates, e.g., that if codeword 01 is formed in bin 4 then we place bits 1,0,1 in bins 3,2,1 respectively. A complete encoder description also requires identification of the probability region over which each bin should be used, or a rule for mapping input bits to bins. We omit this detail to simplify the discussion.

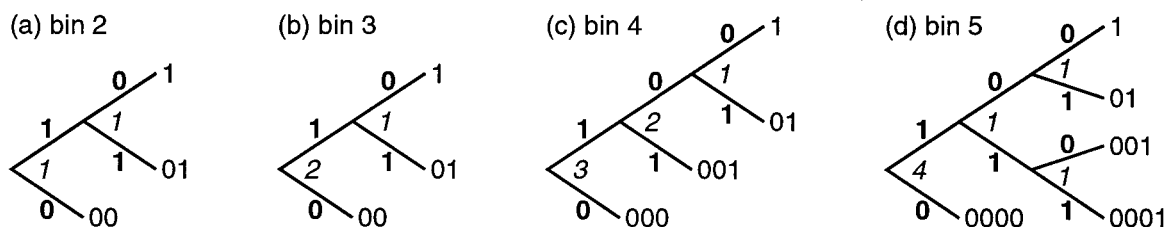


Figure 3: A possible design for a five bin encoder. The first bin is uncoded, hence no tree is shown. Bin indices are shown in italics, output bits in boldface. The input bits are the codewords shown as terminal nodes of each tree.

### 1.3 Relation to Interleaved Entropy Coding

An important special case of the entropy coder arises when all output bits generated from each tree in the encoder are mapped to the uncoded bin. Imposing this restriction reduces encoding complexity, and the encoder amounts to interleaving several separate variable-to-variable length binary codes. This technique was first suggested

in [5], which used interleaved Golomb codes for compression. Howard [4] gives a more thorough analysis of interleaved entropy coding.

By increasing the number and complexity of the variable-to-variable length codes, it's clear that we can make asymptotic redundancy arbitrarily small. With the additional flexibility of the technique presented here, a given redundancy target should be achievable with fewer and/or simpler codes.

## 2 Encoding and Decoding

### 2.1 Decoder Operation

We now describe the encoder and decoder operation, beginning with decoding since it determines the encoding procedure. It is convenient to think of each bin in the decoder as containing a list of bits. To decode, initially we place all of the encoded bits in the first (uncoded) bin, and all other bins are empty. At any time, each nonempty bin (with the exception of the uncoded bin) will contain a single codeword or a suffix of a codeword. Decoding the next source bit amounts to taking the next bit from the bin to which the source bit was assigned. If this bin is empty, we first reconstruct the codeword in that bin by taking bits from other bins as needed.

Software decoding uses two recursive procedures, `GetBit` and `GetCodeword`. `GetBit` simply takes the next available bit from the indicated bin. If the bin is empty then it first calls `GetCodeword`. Given an empty bin, `GetCodeword` determines which codeword must have occupied the bin by taking bits from other bins (via `GetBit`), then places that codeword in the bin. The `GetCodeword` procedure is similar to Huffman decoding, except that at each step we take the next bit from the appropriate bin, not (necessarily) from the encoded bit stream.

To decode the  $i$ th bit, let `binindex` equal the index of the bin to which the  $i$ th bit would have been assigned. This assignment may be a function of any previously decoded bits and any other information available to both the encoder and decoder. Then the  $i$ th decoded bit is equal to `GetBit(binindex)`.

### 2.2 Encoder Operation

To ensure that decoding is possible, we must pay careful attention to the order in which bits are processed by the encoder. Processing bits in the correct order is not straightforward.

One encoding method that produces encoded bits in the appropriate order involves maintaining a linked list of bit values. Each record in the list stores the bit value and the index of the bin that contains the bit. Initially the list contains the entire input sequence in order of arrival. When a codeword is processed, we delete the bits that formed the codeword and insert the resulting output bits in the list at the location of the first bit in the codeword.

To determine the order in which the codewords are processed, perhaps the conceptually simplest method (though in our experience not the fastest) is at each step

to identify the nonempty bin with the highest index. We take bits (in order) from this bin until we have formed a codeword, appending flush bits if needed to complete the final codeword of the bin.

For example, suppose the linked list for the encoder of Figure 3 is as shown in the left half of Figure 4. Bin 4 is the highest indexed nonempty bin, so we search through the list for bits in bin 4 until we form the codeword 01. This codeword produces output bits 1,0,1, in bins 3,2,1 respectively (see Figure 3(c)), so these records are inserted in the linked list as shown in the right half of Figure 4.

When all bits are in the first bin, the encoder output consists of these bits taken in order.

To manage long input sequences with limited memory, we can partition the input sequence into blocks of known size and encode each block separately<sup>3</sup>.

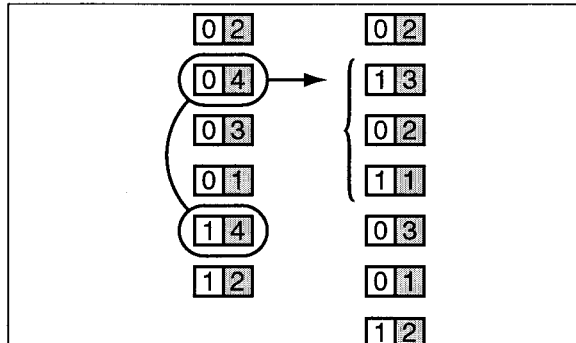


Figure 4: One step of encoding in software using the encoder of Figure 3. In each pair, the left (unshaded) box indicates bit value; the right (shaded) box shows bin index.

### 3 Estimating Rate

We would like to quantify the performance of a given encoder design. One metric we can estimate is the rate (the expected number of output bits per input bit) when the input to the encoder is an independent and identically distributed stream of bits into bin  $j$ , each bit having probability of zero equal to  $p$ . We denote this quantity by  $R_j(p)$ .

Since bins produce output bits that are placed in other bins, estimates of  $R_j(p)$  generally rely on rate estimates for other bins. If bin  $\ell$  has as input  $\lambda_1$  bits with probability  $q_1$  and  $\lambda_2$  bits with probability  $q_2$ , the resulting contribution to rate might be approximated as

1.  $\lambda_1 R_\ell(q_1) + \lambda_2 R_\ell(q_2)$ , or
2.  $(\lambda_1 + \lambda_2) R_\ell \left( \frac{\lambda_1 q_1 + \lambda_2 q_2}{\lambda_1 + \lambda_2} \right)$ .

The first approximation would tend to be more accurate when long runs of bits in bin  $\ell$  have the same probability, and the second would be more accurate if the two types of bits are well mixed. The first tends to be optimistic when  $R_\ell(p)$  is convex  $\cap$ , the second tends to be optimistic when  $R_\ell(p)$  is convex  $\cup$ .

In this section we describe two recursive techniques for estimating  $R_j(p)$  based on the above approximations. Both techniques usually give quite good estimates. The rate estimates produced are asymptotic as the input sequence length becomes large, i.e., the cost of bits used to flush the encoder is not included.

<sup>3</sup>We have developed a more efficient alternative technique but it is beyond the scope of this paper.

The rate estimation techniques do not give exact results because the rate functions are in general nonlinear, and because bits arriving in each bin may not be independent. This dependence arises because encoding a single codeword may result in multiple output bits being placed in the same bin.

### 3.1 First Method for Rate Estimation

We can estimate  $R_j(p)$  recursively using the estimates for  $R_1(p), R_2(p), \dots, R_{j-1}(p)$ . If for each input bit non-terminal node  $k$  of the tree for bin  $j$  produces  $\eta_k(p)$  expected bits, each with probability of zero  $q_k(p)$ , then we use the estimate

$$R_j(p) = \begin{cases} \sum_k \eta_k(p) R_{B(k)}(q_k(p)), & j > 1 \\ 1, & j = 1 \end{cases} \quad (1)$$

where  $B(k)$  is the output bin index for the  $k$ th node in the tree, and the sum is over all non-terminal nodes in the tree.

For example, using (1) to estimate  $R_4(p)$  for the encoder of Figure 3 gives

$$R_4(p) = (1 - p) \left( 4 - \frac{1 + p}{2 + 2p + p^2} + \frac{1}{1 - p^{12}} \right).$$

### 3.2 Second Method for Rate Estimation

In the second technique for estimating  $R_j(p)$ , for each bin we produce a list of  $(\lambda, q)$  pairs. Each pair in the list represents an expected number of bits  $\lambda$  and corresponding probability of zero  $q$  arising from the output of some higher indexed bin or from the source. Initially each list is empty except the list for bin  $j$ , which contains the pair  $(1, p)$ .

At each step, if the list for bin  $\ell$  (initially  $\ell = j$ ) contains pairs  $(\lambda_1, q_1), (\lambda_2, q_2), \dots, (\lambda_m, q_m)$ , we compute the total expected number of bits in the bin

$$\Lambda_\ell = \sum_k \lambda_k$$

and the average probability of a zero in the bin

$$Q_\ell = \frac{1}{\Lambda_\ell} \sum_k q_k \lambda_k.$$

Treating the input to bin  $\ell$  as  $\Lambda_\ell$  bits, each with probability of zero  $Q_\ell$ , we compute the resulting pairs  $(\lambda'_k, q'_k)$  at each non-terminal node in the tree and append  $(\lambda'_k, q'_k)$  to the list for the bin to which the output bit associated with node  $k$  is mapped.

We repeat this procedure, continuing to the first bin. Finally, our estimate of  $R_j(p)$  is equal to the total expected number of bits in the first bin,  $\Lambda_1$ .

Using this technique to estimate  $R_4(p)$  for the encoder of Figure 3 gives

$$R_4(p) = (1 - p) \left( 3 + \frac{(2 - p^6)^2}{(3 - p^2 + p^3 - p^5)(1 - p^6)} \right).$$

For  $p \in (1/2, 1)$ , the method of (1) gives a slightly higher rate estimate in this example.

Variations of these techniques can be used to accurately estimate the rate obtained for a source that produces bits with varying (but known) distributions.

## 4 Code Design

In this section we illustrate a procedure to design a coder. We begin with a redundancy target  $\Delta$  which is the maximum allowed redundancy (in bits per source bit) and a set of candidate trees to be used in the encoder. In this context each tree does not include assignments of bin indices to non-terminal nodes or output bit labels to branches. These assignments will be made as part of the design procedure.

In addition to the family of trees illustrated in Figure 2, there are many other useful trees we can use to design good codes. For example, Figure 5 shows the five useful trees with four terminal nodes.

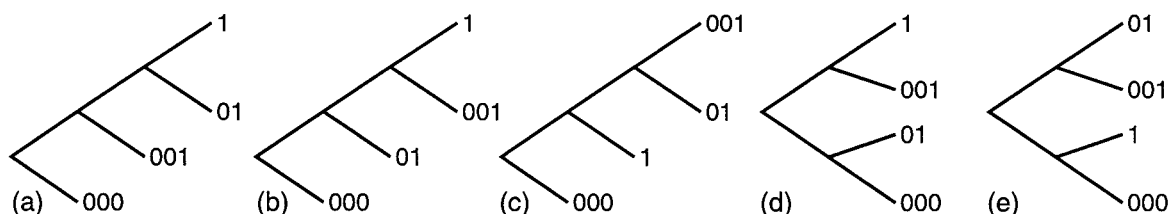


Figure 5: The set of useful trees with four terminal nodes.

We select a trees for each bin in order of increasing bin index. When designs for bins  $1, 2, \dots, j-1$  have been completed, designing the  $j$ th bin amounts to selecting a tree for the bin, assigning bin indices to non-terminal nodes and output bit labels to branches, and calculating  $z_{j-1}$ , the probability value where we switch from bin  $j-1$  to bin  $j$ . (Of course no design work is required for the first bin since it is uncoded and  $z_0 = 1/2$ .) For example, Figure 6 shows a case where the encoder has been designed for the first three bins, and our redundancy target  $\Delta$  is met when  $p$  is less than some value  $p^*$ . Thus we know that  $z_3 \leq p^*$ , and we need to specify the tree to use for the fourth bin of the encoder.

To do this, we can take from our set of candidate trees any tree that is useful at  $p^*$  and assign branch and non-terminal node labels based on this probability value. That is, we calculate the branch probability for each non-terminal node in the tree and label the branches so that a zero output is more likely than a one at each node. Then, at each non-terminal node, if a zero output bit occurs with probability  $q$ , we map this bit to the bin with index  $\ell$  such that  $q \in [z_{\ell-1}, z_\ell]$ .

This construction maps output bits to bins in regions where the redundancy is less than the target  $\Delta$ , and it can be shown that

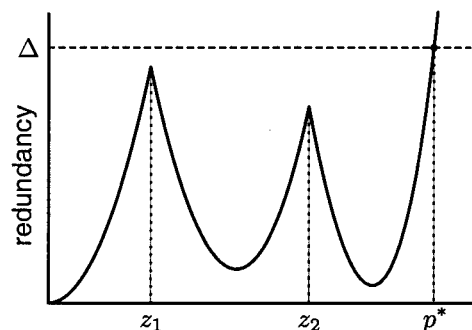


Figure 6: Redundancy of an encoder after designing the first three bins.



(to the extent that (1) is accurate) the redundancy at probability  $p^*$  is strictly less than  $\Delta$ . This follows in part from the following lemma (proof omitted):

**Lemma** *If a tree is useful at  $p^*$ , then the expected number of output bits is less than the expected number of input bits, or equivalently,*

$$\sum_k \lambda_k(p^*) < 1$$

(where  $\lambda_k(p^*)$  is the expected number of output bits generated at node  $k$  and the sum is over all non-terminal nodes  $k$ ).

Thus, the tree we have selected for the new bin produces redundancy less than  $\Delta$  at probability  $p^*$ . Since the rate functions for each bin are continuous, we have extended the range where the encoder meets the redundancy target.

We can also try assigning branch labels, and even selecting a tree, based on some probability target value larger than  $p^*$ . This alternative generally produces larger redundancy at  $p^*$ , but frequently meets the redundancy target  $\Delta$  at  $p^*$  and may extend further the range over which bin  $j$  is used, which can help to reduce the total number of bins used in the encoder.

## 5 Results and Conclusion

Figure 7 shows the redundancy of some entropy coders designed using the technique described in Section 4. The trees in these codes have an average of about 6 terminal nodes.

Figure 8 shows the estimated and measured redundancy of a 23-bin coder which uses only the family of trees illustrated in Figure 2. The estimated and measured redundancy are clearly in close agreement, and the two redundancy estimates obtained using the methods of Section 3 are indistinguishable at this scale.

Decoding speed was tested for a 10-bin coder. For this test, a sequence of probabilities was generated from a uniform distribution on  $[0, 1]$ , and random bits were generated according to these values. We are not concerned with modeling the probabilities here, so bin assignments were computed outside of the timing loop as the optimal assignment given the bit probability; in this way we measured only the speed (and efficiency) of the actual coding.

This coder was tested against the “shift/add” binary arithmetic coder from [6] with parameters  $b = 16$ ,  $f = 12$ . The arithmetic coder was modified to be similarly isolated from the modeling; bit probabilities were supplied in a form convenient to the coder.

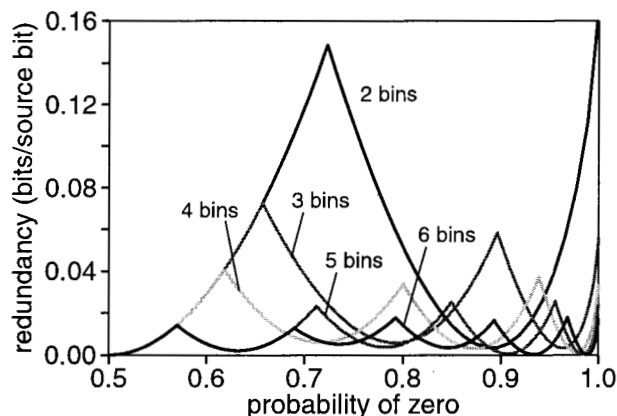


Figure 7: Estimated redundancy of some encoders that have a small number of bins, computed using the method of Section 3.2.

The measured decoding speed for our 10-bin coder was about 2.7 Mbit/sec, compared to 2.0 Mbit/sec for the arithmetic coder. The test used a Sun Ultra Enterprise with a 167 MHz UltraSPARC processor. In this test the redundancy of both coders was very low: 0.0032 bits/source bit for our coder and 0.0019 bits/source bit for the arithmetic coder.

It should be noted that our encoder was about 8 times slower than our decoder. In contrast to the decoding, we made no effort to optimize the encoding; regardless, encoding does appear to be inherently slower than decoding.

Other arithmetic coding variations (e.g. [7]) may be faster than the arithmetic coder we tested. However, our results suggest that, at this early stage of development, our technique offers decoding speeds competitive with those of arithmetic coding.

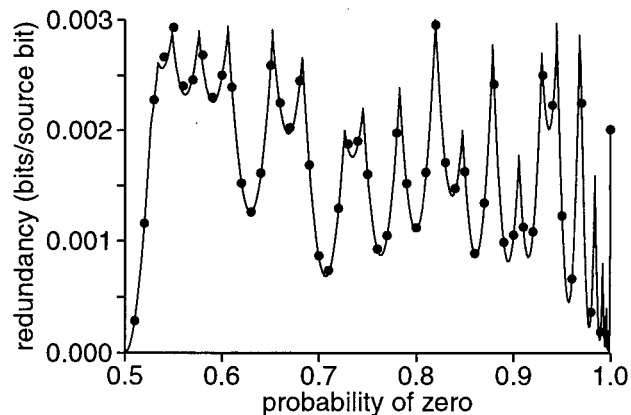


Figure 8: Estimated (solid curve) and measured (individual points) redundancy for a 23-bin coder which has a designed maximum redundancy of 0.003 bits/source bit.

## References

- [1] S. W. Golomb, "Run-Length Encodings," *IEEE Transactions on Information Theory*, vol. IT-12, no. 3, pp. 399-401, July, 1966.
- [2] J. Rissanen and G. G. Langdon, "Arithmetic Coding," *IBM Journal of Research and Development*, vol. 23, no. 2, pp. 149-162, March, 1979.
- [3] I. H. Witten, R. M. Neal, J. G. Cleary, "Arithmetic Coding for Data Compression," *Communications of the ACM*, vol. 30, no. 6, pp. 520-540, June, 1987.
- [4] P. G. Howard, "Interleaving Entropy Codes," *Proc. Compression and Complexity of Sequences 1997*, pp. 45-55, 1998.
- [5] F. Ono, S. Kino, M. Yoshida, and T. Kimura, "Bi-Level Image Coding with MEL-CODE — Comparison of Block Type Code and Arithmetic Type Code," *Proc. IEEE Global Telecommunications Conference (GLOBECOM '89)*, pp. 0255-0260, Nov. 1989.
- [6] A. Moffat, R. M. Neal, and I. H. Witten, "Arithmetic Coding Revisited," *ACM Transactions on Information Systems*, vol. 16, no. 3, pp. 256-294, July, 1998.
- [7] W. B. Pennebaker, J. L. Mitchell, G. G. Langdon, and R. B. Arps, "An Overview of the Basic Principles of the Q-Coder Adaptive Binary Arithmetic Coder," *IBM Journal of Research and Development*, vol. 32, no. 6, pp. 717-726, November, 1988.